

# 1 Inductive Bias

less for transformer

## 2 Weight Reuse and Sharing

1. FC
2. Conv: Sharing in space. Differs when rotating. so rand-rotating req
3. Recurrent: Sharing in time. "Hidden state"

## 3 transformer

### 3.1 tokenize

token set? not exhaustive  
character set, not enough  
subword :bert: word piece, by pair encoding

### 3.2 encoder

Position encoding embedding in context: transformer encoding.

### 3.3 decoder

masked: cannot foresee. Query: from encoder to decoder. project

### 3.4 add&normalize

prenorm (we get it, stable) / postnorm (warmup)

## 4 VIT: Vision transformer

Patch, linear projection (conv, lol)  
class token -> transformer encoder -> MLP head (prenorm)

## 5 attention

$B \times C \times H \times W \rightarrow B \times C \times N \rightarrow B \times N \times C \rightarrow B \times N \times 3C$  (qkv)  $\rightarrow 3 \times B \times N \times C \times D$   
 $B \times H \times N \times D \rightarrow B \times H \times N \times N \rightarrow \text{scale } \frac{1}{\sqrt{h}}$